

# Exploiting the Power of DSM and SMP Clusters for Parallel CFD

Paul R. Woodward, Steven E. Anderson, David H. Porter, Dennis Dingo, Igor Sytine,  
Thomas Ruwart, Michael Jacobs

Laboratory of Computational Science & Engineering, University of Minnesota  
and

R. H. Cohen, B. C. Curtis, W. P. Dannevik, A. M. Dimitis, D. E. Eliason, A. A. Mirin

Lawrence Livermore National Laboratory  
and

Karl-Heinz Winkler and Stephen Hodson  
Los Alamos National Laboratory

## Introduction:

For over a decade, we have been pursuing every available opportunity to use powerful new computing hardware in order to enhance the faithfulness to nature's complexity of our simulations of turbulent fluid flow. The latest challenge in this long process is posed by the powerful new clusters of shared memory multiprocessors (SMPs). Powerful computing platforms of this new type can be found today at Livermore, Los Alamos, NCSA, San Diego, the Minnesota Supercomputing Institute, and elsewhere. If we disregard the distinction between DSM (distributed shared memory) and SMP (symmetric multiprocessor) machines, since it has only a minor impact on programming requirements, we can view the emergence of these computing systems at the largest computing centers in the United States as a welcome convergence of what was heretofore a rather confusing U. S. high-end computer industry.

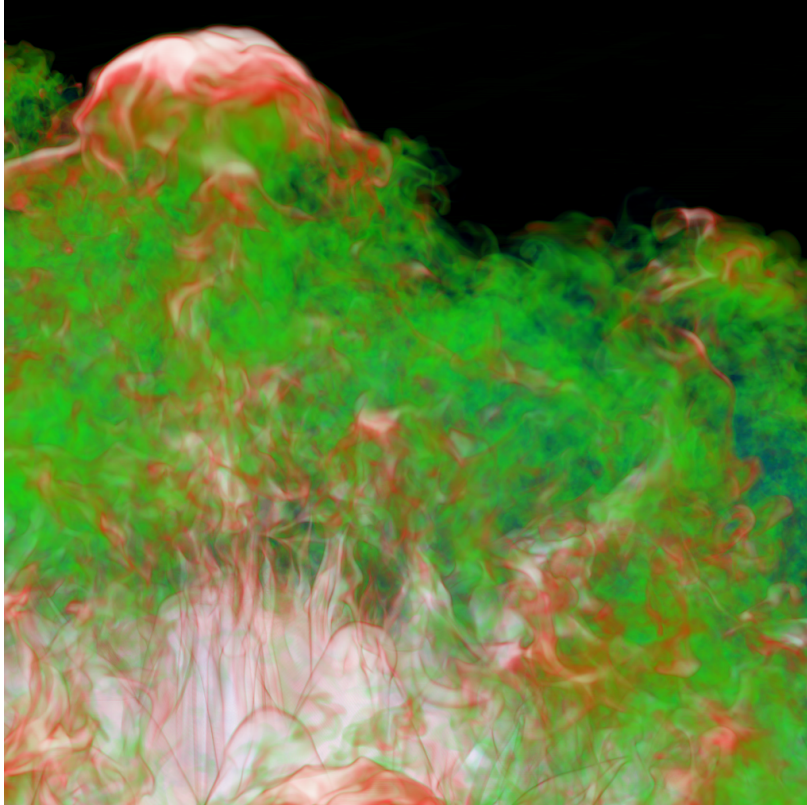
We will report below some of the new 3-D simulations of compressible turbulent fluid flow that we have performed with our PPM gas dynamics code on these DSM cluster platforms. These include simulations of convection in the envelopes of pulsating and rapidly rotating model stars, of chunks of homogeneous, compressible turbulence, and of the unstable acceleration of fluid interfaces by shocks. In order to carry out the most aggressive of these new calculations, we have completely restructured our PPM code to conform to the needs of the new DSM cluster architecture. To do so, we have implemented in software an extremely coarse-grained shared memory model over the cluster. In the most advanced of our implementations, recently demonstrated at NCSA, we have used a hierarchical shared memory programming model to adapt on time scales of seconds to the constantly changing user loads on the DSM machines in the cluster. This allowed us to carry out a billion-cell turbulence calculation as a background job at NCSA, without any sacrifice in performance. We will outline our approach to hierarchical shared memory programming, using this new code and the NCSA demonstration as an example of what we feel is a far more general approach to parallel CFD.

## Finally Resolving Turbulent Flows:

The power of DSM clusters has enabled us during the last two years finally to resolve a reasonable range of turbulent scales within the context of an interesting larger flow. The largest such calculation that achieves this goal is a simulation of the Richtmyer-Meshkov instability of a shock-accelerated fluid interface. This simulation was performed at Livermore with the simplified version, sPPM, of our PPM code running on a grid of 8 billion uniform cells. Because this sPPM code version was used as the ASCI Blue platform procurement's official SMP cluster benchmark, we had the expert assistance of Steve White and his colleagues at IBM. They handled the details of communication between the two (of three) 488-node (1952-processor) "boxes" in Livermore's IBM SP system for this special code implementation. Steve White also collaborated with IBM's compiler group to produce outstanding performance for this code on the entire 5856-processor system. By artfully counting the floating point operations that can, legally, be associated with the reciprocal and square root functions, the aggregate performance was revealed to be in excess of 1 Tflop/s sustained.

The Richtmyer-Meshkov instability results when a shock accelerates an interface separating two

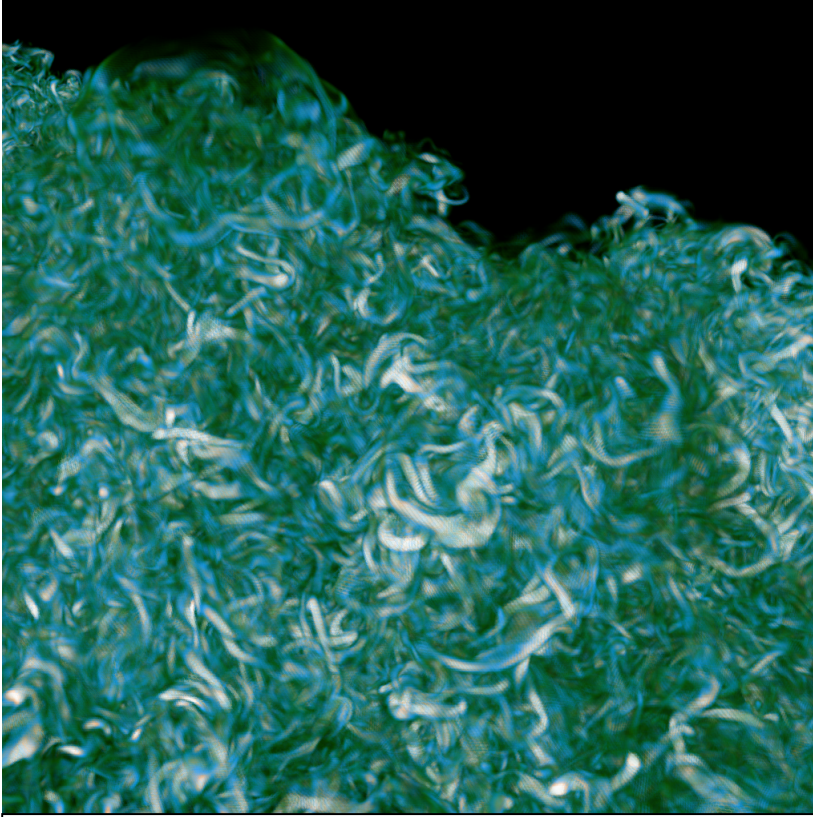
fluids of different densities. Because the shock travels more slowly in the denser medium, perturbations in the interface location cause sections of the shock to lag slightly behind the shock's average position. The oblique sections of the shock which therefore develop set up transverse motions in the compressed gas behind it which are focused on the regions behind the sections of the shock which lag. These focused motions increase the pressure in these regions, and hence the shock strength and propagation speed there. As a result, the shock straightens out, but it leaves behind transverse motions which are converted to longitudinal ones due to the jump in the sound speed at the fluid interface, where the denser material squirts out into the less dense one. These longitudinal motions are not further accelerated, but they are not damped either. Shear at the interface subsequently leads to the development of turbulence and (in the absence of surface tension) to a mixing of the two fluids on small scales.



*Figure 1A. A perspective volume rendering of the distribution of entropy near a small section of the unstable interface near the end of the sPPM simulation of the Richtmyer-Meshkov instability. The entropy of the shocked denser gas is shown as white, while that of the shocked, more diffuse gas is transparent. The region of turbulent mixing is in the green region of this “forest of broccoli.”*

We initialized our simulation of this instability with a planar Mach 1.5 shock approaching through the less dense gas ( $\gamma = 1.3$ ) a perturbed interface with a gas 4.88 times denser ( $\gamma = 1.3$ ). The perturbation of the interface was  $0.01 [-\cos(2\pi x)\cos(2\pi y) + |\sin(10\pi x)\sin(10\pi y)|]$ . We initialized the fluid interface as a smooth transition spread over a width of 5 grid cells. This initialization greatly reduced the amplitude of the high frequency signals that are unavoidable in any grid-based method. The sPPM method of capturing and advecting fluid interfaces forces smearing of these transitions over about 2 grid cells and resists, through its inherent numerical diffusion, development of very short wavelength perturbations. By setting up the initial interface so smoothly, we assured that after its shock compression it would contain only short wavelength perturbations that the sPPM scheme was designed to handle, and that none of these perturbations would be mistaken by the scheme for real signals. Nevertheless, the problem is physically unstable, so one cannot be too

careful in interpreting the results. We are guided in our interpretation by a series of coarser grid simulations of the same problem and by the data shown below. The main point we would like to stress here is that on this grid of 8 billion cells, we are finally at a point where we are able to resolve in this single computation both the primary, long wavelength behavior of the Richtmyer-Meshkov instability and also the secondary, short wavelength behavior of the turbulence that grows out of the shear which this instability produces.

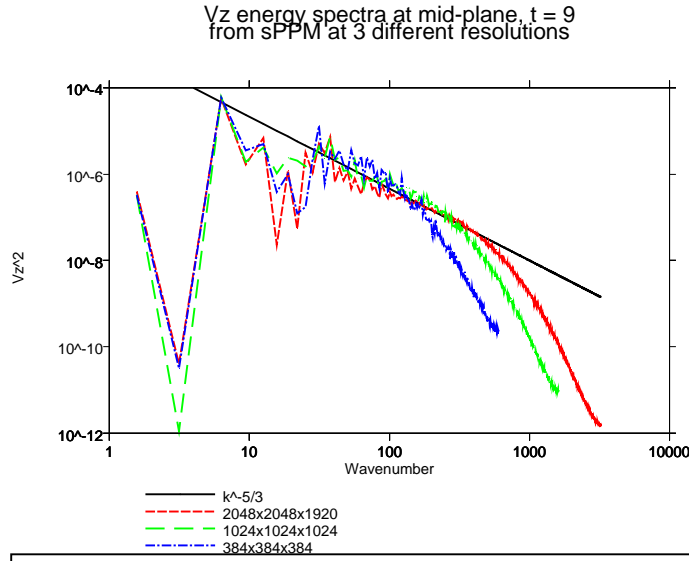


*Figure 1B. A perspective volume rendering of the same region near the unstable interface as is shown in the previous figure. Here the magnitude of the vorticity is visualized. This twisted collection of vortex tubes is characteristic of fully developed turbulence, a conclusion that is supported by the velocity power spectra.*

This sPPM simulation was carried out on a brand new computing system, which therefore did not possess all its planned support systems of disks, archival storage, and visualization hardware. As a result, we were able to archive only 10 full-information snap shot files, each compressed by a factor of 2 to represent each number in only 16 bits. Each of these files was 84 GB in size, so this is still nearly a TB of information. However, based upon preliminary runs of this same problem at lower grid resolution, we determined a single variable, related to the entropy of the gas, and a scaling of this variable to 255 intensity levels, that we wished to save in more complete form. Each snap shot for this single entropy variable was only 8 GB, which made it possible to archive 274 such snap shots during the course of the simulation. We may thus ask any question about the dynamics of the gas entropy in this run and have a good hope of an answer, but questions about other fluid dynamic state variables must be

restricted to their behavior as shown only in the 10 larger snap shots which we were able to preserve. Such constraints are typical of computations carried out on first-of-a-kind computing systems like this one.

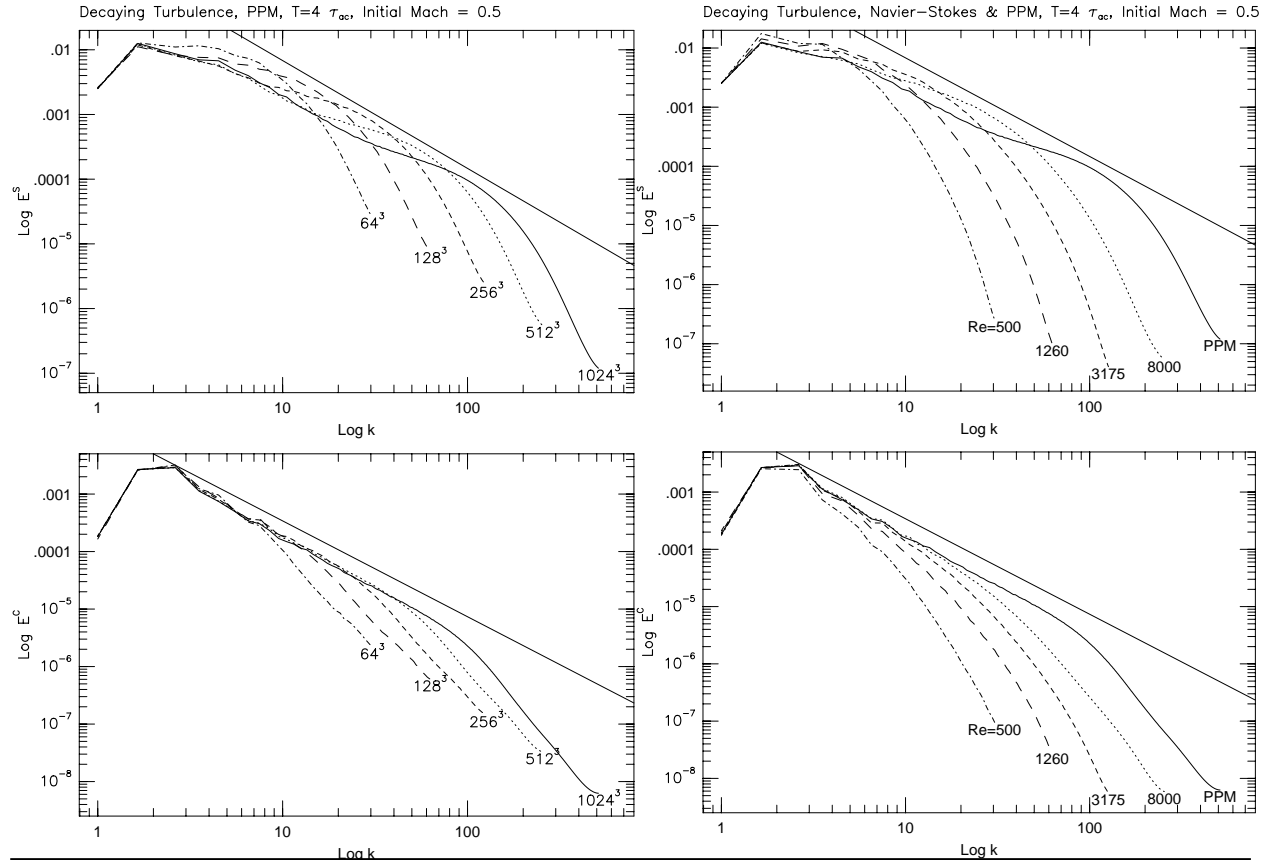
In Figures 1A and 1B above, the entropy and the magnitude of vorticity are visualized in the same small region near the center of the unstable interface near the end of the simulation. From these figures we can see that this simulation has indeed captured both the macroscopic and the microscopic scales of the Richtmyer-Meshkov instability and its secondary Kelvin-Helmholtz instabilities. We can make this statement quantitative by looking at the velocity power spectrum within a plane slicing through the unstable layer. Such a power spectrum is shown in Figure 2. Actually, in this figure three such spectra are compared. The three spectra come from simulations on progressively finer grids of  $384^3$ ,  $1024^3$ , and  $2048^2 \times 1920$  cells. The low frequency parts of these spectra reflect the initial perturbations and the harmonic modes that their nonlinear interactions have produced. To the right of the higher frequency initial perturbation, a short section of a turbulent inertial range, with a Kolmogorov power-law slope of  $k^{-5/3}$  can be identified. The coarsest grid does not reveal such an inertial range, but both of the finer grids show this behavior. Still further toward high frequencies, in the near dissipation regime, the power spectra flatten somewhat, developing power-law behavior more like  $k^{-1}$ , characteristic of the Fourier transform of very long vortex tubes (as are visible in the image of the vorticity). On the finest grid, we



have an inertial range of the turbulence of about a factor of 3 in extent, with the spectrum flattening toward  $k^{-1}$  in the near dissipation range and then steepening greatly in the dissipation range (for wavelengths of about 8 grid cells or less). Our experience with simulations of homogeneous turbulence (see below) indicates that the indirect effects of the dissipation of the turbulent motions by the numerical viscosity of the sPPM scheme extend no further toward long wavelengths than the short section of  $k^{-1}$  slope. Therefore, presumably, this simulation properly incorporates the physical effects of the interactions of the original perturbation scales and their first few harmonics with the smaller scales of the turbulence that the instability has generated.

#### Concentrating all the DSM Cluster Computational Power on a Single Chunk of Turbulence:

Using the Silicon Graphics DSM cluster at Los Alamos in the spring and summer of 1997, we were able to simulate a small chunk of homogeneous, compressible turbulence. With the PPM gas dynamics code running with periodic boundary conditions on a uniform grid of  $1024^3$  cells, we achieved a fully developed turbulent flow containing a Kolmogorov inertial range of about a factor of 8 in wavenumber.



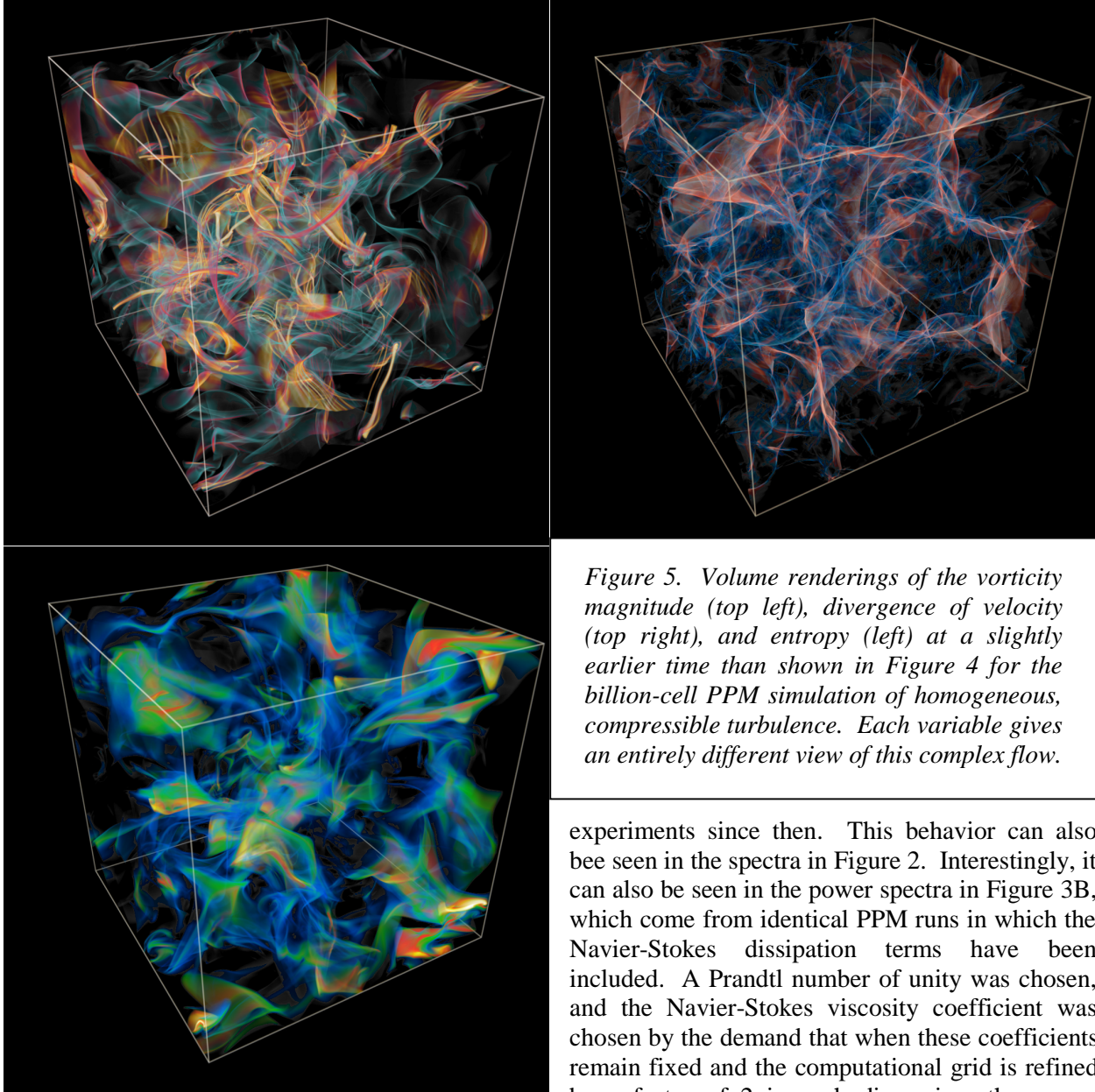
Exploiting the Power of DSM and SMP Clusters for Parallel CFD





*Figure 4. The distribution of the magnitude of vorticity is visualized in this volume rendering of a region of the billion-cell simulation of homogeneous, Mach  $\frac{1}{2}$  turbulence. Near the center of the figure, a group of thin vortex tubes are spiraling around each other as a whole section of a vortex sheet rolls up. At the left, a series of roughly parallel vortex tubes is intensifying out of another vortex sheet. This visualization captures the developing turbulent flow in the midst of its transition from a state dominated by vortex sheets to one in which such sheets can no longer be identified within the dense tangle of spaghetti-like vortex tubes.*

This flow was initiated with a statistical sample of long wavelength velocity disturbances centered fairly narrowly on a wavelength of half the periodic scale. The density and pressure were constant in the initial state, and the amplitude of the initial velocity disturbances gave an rms velocity perturbation of half the sound speed. The velocity power spectrum is shown in figure 3A at a point during the decay of this flow when the turbulence has become fully established. This spectrum is compared in the figure with those of identical runs with the PPM code on grids of progressively coarser resolution:  $512^3$ ,  $256^3$ ,  $128^3$ , and  $64^3$ . The velocity power spectra are shown in this figure for both the incompressible and compressible components of the velocity field (lower and upper panels). The convergence of the power spectra is clear. Each successively finer grid leaves the long wavelength behavior invariant while adding another section to the spectrum at high wavenumbers. At the right-hand end of each power spectrum is of course a short region of strong damping due to the numerical dissipation of the PPM Euler scheme. For the solenoidal power spectra, this picture is complicated by the presence in the near dissipation range of a section of flatter slope, about  $k^{-1}$ , extending over roughly a factor of 4 in wavenumber. This behavior was noted first in our numerical simulations several years ago, and it has been confirmed in other simulations as well as in

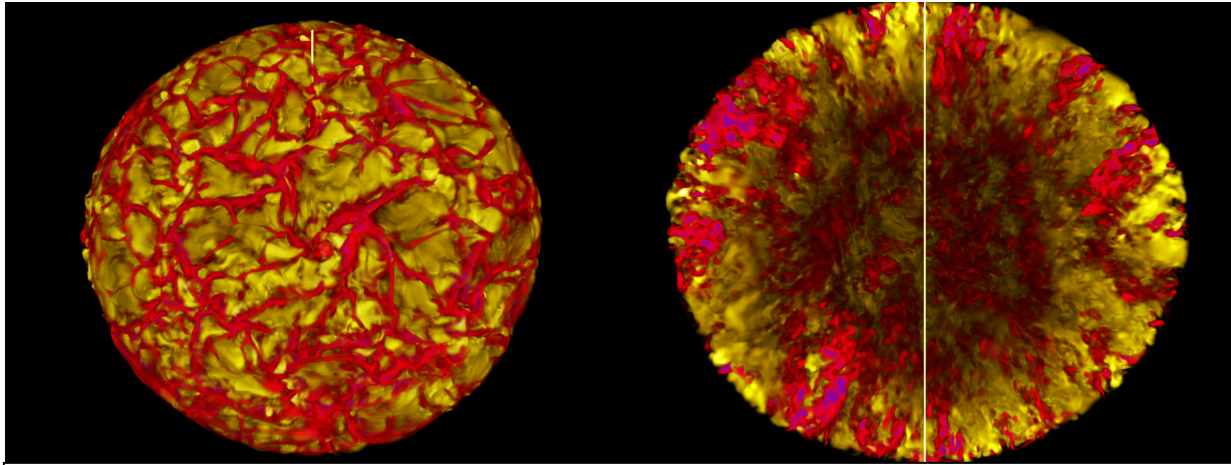


*Figure 5. Volume renderings of the vorticity magnitude (top left), divergence of velocity (top right), and entropy (left) at a slightly earlier time than shown in Figure 4 for the billion-cell PPM simulation of homogeneous, compressible turbulence. Each variable gives an entirely different view of this complex flow.*

spectra remain unchanged. This demand of convergence of the power spectra can of course not be applied in the Euler case, since a grid refinement reduces the effective viscosity of the flow by a factor of 8 (for PPM). The Navier-Stokes runs were carried only to a grid of  $512^3$ , at which resolution there is still no Kolmogorov inertial range at all for this flow. Nevertheless, the trend of these curves — for grids of  $64^3$ ,  $128^3$ ,  $256^3$ , and  $512^3$  — indicates that with a grid of about  $4096^3$  we would see such a region of the spectrum with about the same extent in wavenumber as the Euler approach gives us with only  $1024^3$  cells.

The billion-cell Mach  $\frac{1}{2}$  turbulence calculation gives an enchanting, detailed view of the development of a turbulent flow. The instability of originally smooth vortex sheets to form systems of vortex tubes and the subsequent braiding of these tubes about each other is clearly seen. Figure 4, above, gives a glimpse of this fascinating process when it is about half-way along, and Figure 5 shows volume renderings of three variables, the vorticity magnitude, the divergence of the velocity, and the entropy, respectively from left to right, a bit earlier in the simulation.



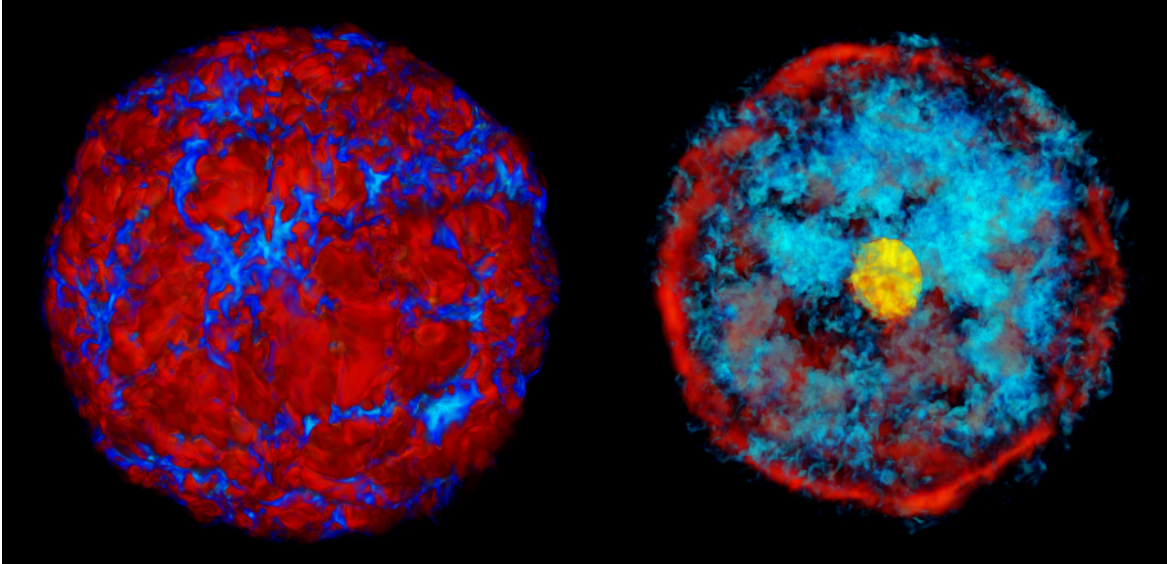


*Figure 6A,B. The distribution of the vertical velocity component is shown in these two volume renderings from a PPM simulation of turbulent compressible convection in a rapidly rotating model star. This calculation was performed on a uniform Cartesian grid of  $512^3$  cells. The yellow regions of the flow are moving upward while the downflows are represented as red and purple. The rotation axis is represented as a white line. The outer view, at the left, shows the cellular convection pattern near the surface of the star. The oblateness of this very rapidly rotating object is more apparent in the view at the right from the equatorial plane. In this volume rendering, the near half of the star has been cut away, and the central half by radius is made transparent. This central region of the star is convectively stable. The eye can detect in the right-hand image a tendency for the lanes of upward- and downward-flowing gas to line up roughly along longitudinal lines. This tendency to form banana-shaped convection cells is countered by the differential rotation of the star, which has in turn resulted from the convective transport of angular momentum. An animation of the convection patterns as viewed from the stellar interior reveals a continual forming, tearing, and reforming of these banana structures.*

#### Simulating Turbulence Driven by Convection, which in turn Interacts with Stellar Rotation or Pulsation:

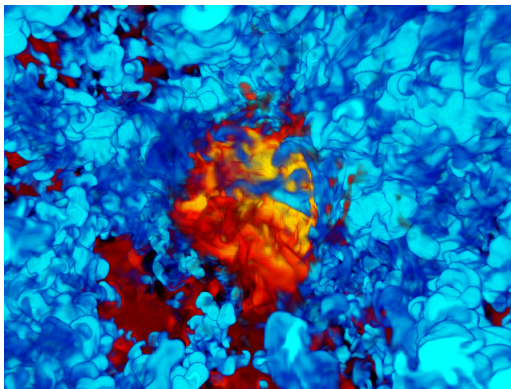
In a number of PPM simulations using NCSA's Silicon Graphics DSM cluster, we have explored the interaction of turbulence with thermal convection, and of this turbulent convection in turn with either rotation or pulsation. This work began with our participation in an NSF-funded Grand Challenge team project, together with the team of Juri Toomre at the University of Colorado, Boulder. The computational power of the DSM cluster made it possible to extend the grand challenge work on small 3-D blocks representative of the solar convection zone near its surface to include the entire model star. Our first simulations, begun in the fall of 1996, used a model star which was chosen to explore the interaction of convection with rotation. We rather arbitrarily chose to make the central half of the star (by radius) stable to convection, while the outer half was strongly convectively unstable. To maximize the coupling to rotation, and to make the model star easier to simulate with our explicit gas dynamics techniques, we made the star very luminous (very strong driving for the convection) and we set it initially rotating uniformly with velocities near the equator comparable to the local sound speed there. This calculation incorporated a treatment of the free surface of the star which used a multifluid algorithm based on the SLIC algorithm of Noh and Woodward, but with gravity serving to make the interface between the stellar gas and the second "fluid" (vacuum) stable. The escape of heat through this surface was treated rather crudely, however, with the time-averaged surface heat flux forced to match the constant rate at which heat was introduced into the stable central region of the model star.

A view of this rotating star from just above its equatorial plane is shown in Figure 6A, and a view



*Figure 7A,B. Two perspective volume renderings of a PPM simulation of the convective envelope of a model giant star. Temperature fluctuations relative to average values on isopressure surfaces are shown, with red and yellow depicting warm and hot temperatures and blue and aqua representing cool and cold temperatures. At the left, the envelope is made relatively opaque, so that the surface pattern of convection cells is visualized. At the right, the gas has been made relatively transparent, so that the interior strong dipolar flow pattern is revealed.*

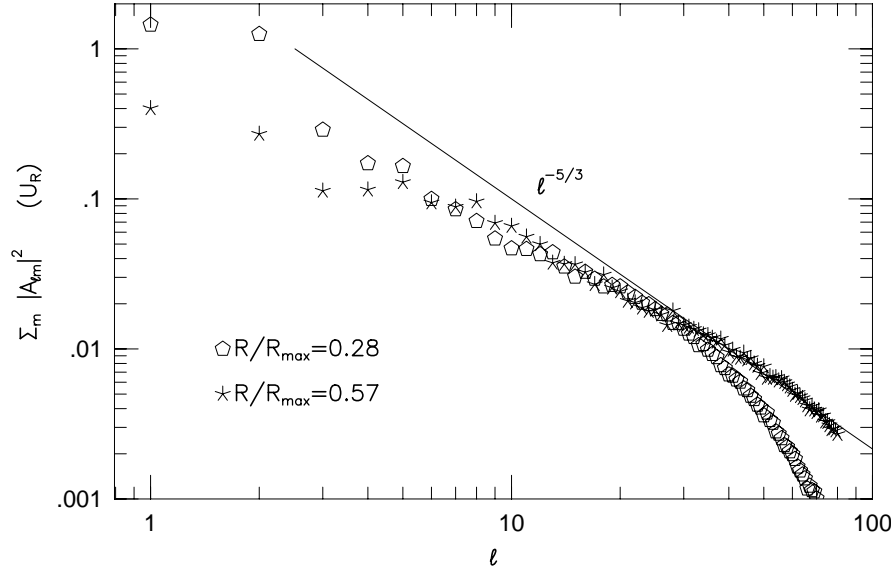
from the equatorial plane with the near half of the star cut away is shown in Figure 6B. In both figures, the rotation axis is represented by a white line. The surface of the star has a cellular convection pattern. The cells are larger than we might expect, but this is due to the relatively coarse  $512^3$  Cartesian grid of this calculation. The convection zone is just 128 grid cells thick in this case, which our earlier simulations of isolated chunks of convection indicate to be adequate but coarse. The convection has transported angular momentum in the gas, so that an easily noticeable differential rotation has developed. Animations of cutaway views of the convection show that the tendency of the convection cells to develop the “banana cell” structure due to the Taylor-Proudman effect is defeated by the strong differential rotation. Such banana structures can often be seen as transients, continually forming, being torn apart by differential rotation, and then reforming from different component elements. In the cutaway view in Figure 6B above, the central half of the star by radius is transparent, since it does not participate in the convection flow. Thus in the central part of that image we are looking at the base of the convection zone from the inside. This calculation treats the surface of the star as a free surface stabilized by gravitational forces. One



*Figure 8. A volume rendering of the convection flow near the stable, hot stellar core, which appears as a red and yellow region partially enveloped by the funnel of cool gas, shown in blue and aqua, striking it from the right. The geometrical convergence of the convection flow near the core eliminates all cellular convection structure in this region, leaving only a very strong, but turbulent, wind from one side.*



of the initial surprises of this simulation was the thin outer belt of gas near the equator that rotates especially rapidly. Examination of the 2 TB of data archived from this run revealed that this gas tends to



*Figure 9. Velocity power spectra from the middle of the convective envelope of our model giant star are here analyzed in terms of spherical harmonic modes. Note the dominance of the dipole mode and the Kolmogorov power-law behavior at higher wavenumbers. The slight flattening of the spectra toward the highest wavenumbers, just before the dissipation range, is a feature we have seen in every compressible, turbulent flow we have simulated.*

collect at the surface in the equatorial region because of its greater buoyancy due to the effect of the greater centrifugal forces exerted on it, which weaken the effective gravity. Such relatively rapidly rotating equatorial outer layers of gas are observed on the sun and on rapidly rotating planets like Jupiter.

In the fall of 1997, we used the same multifluid, 3-D Cartesian PPM code at NCSA to simulate pulsating stars. We again chose a luminous model star to drive very rapid, vigorous convection. In this case, we set up the unperturbed stellar model without rotation but with a very deep convection zone as is found in giant stars. In order to resolve it on the computational mesh, we enlarged the convective-

ly stable, hot stellar core to about 20% of the stellar radius. In a typical giant star, the hot, stable core would have a radius comparable to that of the earth while the convective envelope around it would have a radius comparable to that of the orbit of the earth about the sun. We artificially reduced this tremendous dynamic range in our simulation, keeping the stellar core small but still resolvable. In Figure 7, above, we see two volume renderings of our simulated giant star convective envelope. At the left, we have made the envelope relatively opaque, so that we see the surface features, while on the right we have made the envelope sufficiently transparent that we may see right through it to the hot stellar core within. In both renderings, we show the temperature fluctuations relative to average temperatures on each isopressure surface. Red and yellow represent warm and hot temperatures, while blue and aqua represent cool and cold ones.

Aside from the pulsation of this model giant star envelope over a range of about 20% in radius, the biggest surprise of this calculation was the prominence of global modes of convection. Treatments of convection for evolution calculations for stars of this type assume that the convection can be characterized by mixing length models which involve only local parameters. However, this simulation reveals a propensity for global modes of convection. In fact, animated sequences of images like these reveal a general dipolar flow pattern, with relatively cool gas descending toward the core from the upper right in Figure 7B, becoming heated while passing over the core at about a fifth of the local sound speed, and rising as relatively warm gas to the lower left in the figure. A close-up view of the core region showing this flow appears in Figure 8.

The deep convection in the model giant star envelope that we have been discussing has provided us with a fully developed thermal convection flow with 5 pressure scale heights over which the special conditions of either the upper or the lower boundaries exert no major effects. In this flow the periodic boundary conditions in the “horizontal” directions are not at all artificial, but instead simply express the overall spherical geometry of the problem. Here we have, as in the Richtmyer-Meshkov problem, turbulence driven by a real physical process which is itself simulated in complete detail within the calculation. In contrast to the Richtmyer-Meshkov example, however, this thermal convection flow is statistically steady, except of course for the periodic radial pulsation of the envelope as a whole. This PPM simulation on a  $512^3$  grid was carried out over about 20 pulsation periods, so that the flow has had a good deal of time to relax. Velocity power spectra of the flow near the middle of the envelope (middle values of the radius) are shown in Figure 9. Once again, the longest wavelength modes are characteristic of the driving physical process and there is a turbulent Kolmogorov inertial range, in this case rather short. Again there is a flattening of the power spectra in the near dissipation range followed by a dramatic steepening as the turbulence is dissipated at the highest wavenumbers. In this simulation, not only is the nonlinear interaction of the small scale turbulence with the convection flow treated in detail, but the nonlinear interaction of the dipolar convection flow with the global radial pulsation is accounted for as well. We are beginning a new series of such giant star simulations incorporating more realistic models of the gas equation of state, including gas ionization effects, and improving our treatment of the escape of heat from the stellar surface.

#### The DSM Cluster Programming Challenge:

Today’s high-end computing platforms, DSM and SMP cluster systems like those from Silicon Graphics and IBM which were used to perform the PPM simulations discussed above, combine deep memory hierarchies in both latency and bandwidth with a need for many-hundred-fold to several-thousand-fold parallelism. Users of these systems have had to meet these challenges to efficient parallel program design armed only with minimal system software: Fortran, C, MPI, and support for POSIX threads on a network node. OpenMP is a promising new standard which can be used to generate portable code for a single DSM or SMP machine, but it does not address the cluster as a whole. Programming for dedicated cluster systems would be difficult enough, however DSM cluster systems are usually administered in order to optimize throughput, so that a mix of jobs of different sizes is set running at any given time. It is therefore very difficult to obtain for a single large job even an entire 128-processor machine, let alone a cluster of such machines. In order to run efficiently in this context, we restructured our PPM gas dynamics program so that it could dynamically adjust the number of processors it used on any single DSM machine of the cluster. This allowed the code to coexist with a time varying mix of other jobs. Our code’s task manager continually entered requests into multiple job queues, adding any processors that were made available through this mechanism to the ongoing team. So far we have used this restructured code for problems with a regular structure, adapting to irregular processor loads from other users. However, it will be no additional difficulty to have this code dynamically adjust processor loads as a result of dynamically changing computational loads from regions of our own problem domain. The techniques that we have used to accomplish this code restructuring are outlined below. They are very general and should apply equally well to many other computational problems.

#### Hierarchical Shared Memory:

Our first versions of our PPM gas dynamics code, like the sPPM benchmark code that we wrote for the DoE’s ASCI program, used thread-based shared memory multitasking within each DSM machine and MPI message passing over the DSM cluster. Not only was this hybrid coding technique clumsy, but it also made load balancing over the cluster network difficult. Our present approach extends to the entire cluster the shared memory multitasking approach that we use within each DSM machine. The key to this kind of parallel program is to decompose the work of the program into a sequence of tasks each of which

requires only data from a restricted and compact data context and each of which can be executed from this data context without the need to communicate with any other task. This is a shared memory paradigm; tasks do not communicate directly with each other, but instead they read and write shared memory data structures. If we wish, we can think of the task's data context as including within it "messages" that have been written there by other tasks. This message passing through the intermediary of shared memory is, however, much simpler than message passing directly between ongoing processes. No message receiver ever need know the identity of the message sender, and vice versa. There is also no need for message buffering, since the data structures in shared memory are pre-allocated. It is still a good idea, of course, to write data needed by other tasks as soon as this data is generated and to read data supplied by other tasks at the latest possible moment. Synchronization of this form of message passing is therefore still required, although it is generally much simpler. In our code, we accomplish this synchronization by having each task set a semaphore variable in shared memory indicating when the task has been completed. Tests on these task completion semaphores are performed before each task launch, so that a task once begun knows that it is safe to do all its work without further inquiry. The art of writing such programs is to order the list of tasks very carefully so that at any point in the program a very large number of tasks near that point in the sequence can be executing in parallel.

So far, the strategy of task parallelism outlined above should be familiar. This is the method for constructing parallel programs for SMPs. For DSM machines, there is one further detail. The non-uniform memory access of the DSM architecture forces the programmer to take the trouble of copying into the local memory of the executing processor elements of the task data context that will be overwritten. Any intermediate data generated by the task that will not be written into shared memory for other tasks to read must also be placed in the local memory of the executing processor. In Fortran, this local memory placement is easily accomplished. The task is merely encapsulated in a subroutine (subroutine linkage is a negligible cost for any task with a hope of being efficient), and the data that must be local is dimensioned locally, so that it will be placed by the compiler on the subroutine stack (which is always in the best possible memory). This placement of the task work space in local memory eliminates the phenomenon of "false sharing" and greatly improves performance. Note that it is unnecessary to know which processor will execute the task. The subroutine stack will transparently be put in the right place.

Our program consists of a hierarchical set of task lists. The first set is a list of tasks intended for execution by entire machines of the cluster. These tasks, which we will call DSM tasks, are themselves composed of lists of smaller tasks, which we will call CPU tasks. A CPU task is encapsulated in a subroutine and is written for execution by a single CPU. DSM tasks are encapsulated in subroutines that are explicitly multithreaded for parallel execution. We make each task conform to a template consisting of three steps: (1) reading the task data context into the local, private task work space, (2) operating on the data context, and (3) writing data back into shared data structures. CPU tasks enjoy relatively large data bandwidth to the shared memory of their DSM machines, so they may perform these three task stages sequentially without significant loss of performance. (Machines that do not support this mode of task operation are difficult to sell and therefore are difficult to find.) However, DSM tasks do not enjoy high bandwidth access to the memory of other cluster members (or to shared disks). Therefore, steps 1 and 3 above must be overlapped with step 2. That is, the data for the next DSM task must be prefetched during the execution of the present DSM task. Also, the data produced by this DSM task must be written back to shared data structures during the execution of the next DSM task. This can be accomplished by encapsulating these data transfers in separate DSM tasks, with the obvious constraint that they must be executed by the same DSM that performs the real computational work of the DSM task to which they correspond. In our implementation, we have constructed memory server daemons that run on each DSM machine of the cluster and which asynchronously fulfill requests to "put" and "get" contiguous sections of arrays registered with them as globally accessible. We have coordinated the DSM tasks through a task manager process, which has only a single thread.



### Hierarchical Shared Memory Parallel Implementation of PPM at NCSA:

We restructured the PPM gas dynamics code according to the hierarchical shared memory strategy outlined above. Memory server daemons were created to read and write a fast Fibre Channel network-attached disk system of our own design. Two 128-processor Silicon Graphics Origin-2000 machines at NCSA, interconnected by a single fast Ethernet, shared a common file system on 48 Seagate Fibre Channel disks supplied by LCSE industrial partner MTI. Each machine was connected to all 48 disks via 4 Fibre Channel loops. Each machine was connected to the disks through its own set of ports (the disks were dual ported). Read/write bandwidth from the PPM application from each machine was in excess of 270 MB/s, sustained, even when both machines accessed the disks simultaneously. Control information, such as DSM task completion semaphores, was passed via MPI over the fast Ethernet link.

This restructured PPM code was used to simulate Mach 2 homogeneous, compressible turbulence on a billion-cell ( $1024^3$ ) uniform grid. A typical task for a single CPU was to update for a single 1-D pass a grid pencil of  $4 \times 4 \times 256$  cells. A typical task for a single 128-processor Origin-2000 machine was to update for six 1-D passes, or 2 time steps, a  $256 \times 256 \times 512$  brick of grid cells. The active data context for the job consisted of 32 old and 32 new grid brick records stored on the 864 GB shared Fibre Channel disk system. Each grid brick record of 954 MB consisted of 27 separate records: the brick interior (640 MB), 6 brick face records (27.5 MB each), 12 brick edge records (2.4 MB each), and 8 brick corner records (200 KB each). The active memory context in each participating DSM machine consisted of 5 grid brick records, or about 5 GB out of the 64 GB of DSM memory in each machine.

During each grid brick update, the Origin-2000 was asynchronously prefetching the next grid brick record and writing back the results of the previous grid brick update to 27 different grid brick records on disk. The grid brick record, after being read into DSM memory from disk (in 3.5 sec), was unpacked to form a single, augmented grid brick of  $300 \times 300 \times 556$  cells. This brick was then updated in six 1-D passes, with each consisting of 8192 single-CPU tasks (requiring 2.5 sec with 128 CPUs). In this demonstration run, there were barrier synchronization points at the ends of the 6 passes, but these can be eliminated at the cost of further code complexity. After the 6 passes, the new data was written into a new grid brick record in DSM memory, and this was transferred back to disk (in 3.5 sec).

Figure 10, on the next page, documents about 4 days of continuous PPM computation at NCSA. Two 128-processor Origin-2000 systems were used. Processors obtained by PPM on the first system are represented by the cream colored area in the figure. This system was not always available, due to scheduling of dedicated access for other jobs. Processors obtained by PPM on the second 128-processor system are represented by the blue area in the figure. PPM adjusted its number of processors on each machine at the beginning of each 1-D sweep for each grid brick. When 128 CPUs were in use, this adjustment interval was about 2.5 seconds. Both machines were shared dynamically with other users, and PPM benefited by inserting requests for CPUs in several batch queues, grabbing the CPUs as they became available for this single, large computation. The small departures from full resource utilization that are shown reflect system functions performed by the operators, not any failure of PPM to exploit these opportunities.

### Reevaluating this Demonstration Code for Parameters of the IBM SP System:

Our sPPM benchmark code indicates that, within a couple of per cent, the delivered performance of a PowerPC 604e microprocessor running at 333 MHz is equivalent to that of a MIPS R10K processor running at 195 MHz. Therefore, one SMP of 4 such PowerPC CPUs is about 32 times less powerful than a 128-CPU Origin-2000 for our PPM code. Hence we should reduce the computational labor involved in a single DSM task of the PPM code by about a factor of 32 in order to run well on this platform. We can do this by reducing the brick volume by factor of 8, resulting in a brick of  $128 \times 128 \times 256$  cells, and by performing only a single 1-D sweep rather than 6 per task. This DSM task should now take  $15 \times 32 / 48 = 10$  sec. Because we are performing only a single 1-D sweep, the data reuse in this task is 6 times less than

on the Origin-2000. However, the ratio of “cluster” network bandwidth to DSM processing speed is now  $4 \times 25 / 600 = 1/6$  Bytes/flop. This is 12 times greater than for the Origin-2000 at NCSA. As a result, the PPM code should run on this system even more efficiently. A form of overhead for the parallel code is redundant computation performed in “ghost” cells surrounding each grid brick. The fraction of the computation time devoted to this redundant work would have remained the same as in the NCSA run if we had performed three rather than just one sweep per DSM task. We have thus reduced the redundant computation overhead by a factor of  $(278^2 \times 534 - 256^2 \times 512) / (8 \times (130^2 \times 263 - 128^2 \times 256)) = 3.9$  and the efficiency of the job should therefore be even greater. We note that these projections are merely educated guesses, and they do not incorporate any consideration of problem I/O. They nevertheless suggest that our parallel programming model might be portable to systems with fairly dramatically different parameters than the Silicon Graphics Origin cluster at NCSA.

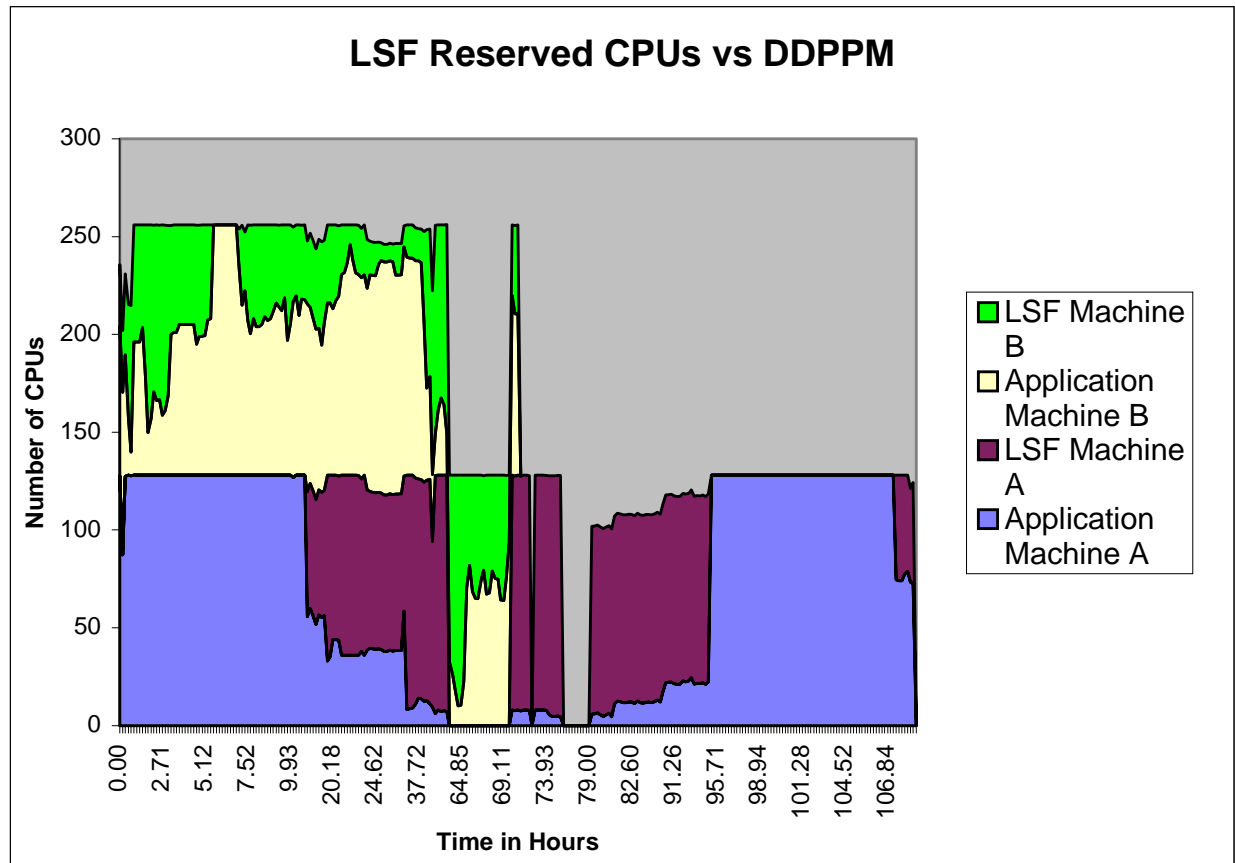


Figure 10. The usage of processors on two 128-processor Silicon Graphics Origin-2000 machines at NCSA is shown here over a 4-day period in the fall of 1998. The light blue and the tan areas of the diagram represent processors dynamically obtained by our PPM code for use on a single, billion-cell simulation of homogeneous, compressible turbulence. The other colors show processors devoted to other users of these systems. About half way through this period, one of the machines was given over to another user for a dedicated run. For a brief interval, both machines were unavailable to us. Brief periods in which processor utilization falls below 100% were caused by use of the machines by the system operators. The minimum response time of this PPM job to changes in resource availability was 2.5 seconds.

### Assumptions that Permit Efficient DSM Cluster Programs:

It is tempting to speculate that perhaps all codes aimed at the simulation of physical systems on grids could be implemented in the above fashion. As a result, we have attempted to abstract those characteristics of our PPM code that we feel are essential in enabling this restructuring. We state them here in the form of assumptions that we believe are necessary for such hierarchical shared memory parallel program execution. First, we assume that a job can be decomposed into a set of tasks that can be executed independently, so long as certain previous tasks are completed at task launch. We further assume that each task can be made to conform to a model, or template, in which:

- 1) possibly remote data is copied into local memory,
- 2) this data is operated upon mightily,
- 3) a few results are written back to possibly remote storage.

We assume that the tasks can be constructed so that, in general, the larger the data context for the task, the larger the amount of potential data reuse. This assumption is necessary to accommodate low cluster bandwidths. To accommodate large cluster latencies we must also assume that the task data contexts can be constructed so that they may be read or written back in only a small number of sequential data transfers. Once in fast local memory, these data contexts can be efficiently reorganized if necessary.

We assume that global barriers (Amdahl's Law) can be avoided by providing greater system resources and/or by minor modifications of the numerical algorithm. Examples of this principle abound. For example, a program may require all processing to stop so that an image of the problem can be written to a restart file on disk. However, if additional system memory is provided, this restart dump can be written asynchronously without impeding the program flow. Another program might require that a global reduction operation be performed after a time step is completed in order to determine the value of the next time step. However, if enough memory is provided to store the previous problem state, we may guess the time step value (the minimum of the previous 25 time steps might be a good guess) and proceed speculatively. In the rare event that we guess badly, the saved system state will permit us to recover. As a final example, we may be performing an implicit calculation that appears to require global information to be assembled in order to update the value of a variable at a single spatial location. By revising the numerical algorithm slightly, we could require up-to-date information only for the local region and use information from the previous time step or iteration for the more distant data. Once again, this would require a commitment of additional system memory to the job.

### Feasibility Requirements:

There are bandwidth requirements for the cluster network, but essentially no meaningful latency requirements. The bandwidth requirements are determined by the demand that any computing resource should be able to execute any task, regardless of the location of its data context. Data prefetching and asynchronous write back are absolutely essential. The task manager needs to have limited intelligence to avoid stupid data movement. It needs to dynamically reorder the task list, permuting elements that are equally or nearly equally qualified candidates for the next task to be launched, taking data location over the network into account. Finally, local memory for various computing resources must be sufficient to accommodate data contexts offering sufficient data reuse, but this is not a new requirement.

### Conclusions:

The new generation of powerful DSM and SMP cluster computers enables simulations of fluid dynamics at sufficient resolution to compute the complex nonlinear interactions of small-scale turbulent motions with a large-scale driving flow. With a new programming model of hierarchical shared memory multitasking, it is possible to exploit these new systems without disrupting the flow of small and medium-sized jobs that makes their existence possible.



### Acknowledgements:

We would like to acknowledge generous support for this work from the Department of Energy, the National Science Foundation, and NASA. Our work on implementing codes efficiently on DSM clusters began in 1993 with support from Silicon Graphics and the Army Research Laboratory and has continued with support from the DoE's ASCI program through a Level-2 project with Los Alamos, subcontract B33700016-3Y, completed in 1998, and with support from NSF through a MetaCenter Regional Alliance grant, ASC-9523480, and support from the PACI program through subcontracts from NCSA. Support for our continued numerical algorithm development and the packaging of our algorithms in library form, PPMLIB, has come from the DoE's Office of Energy Research through grants DE-FG02-87ER25035 and DE-FG02-94ER25207, respectively. Support for our investigations of turbulent compressible convection has come from an NSF Grand Challenge Application Group award, ASC-9217394, through a subcontract from the University of Colorado, and more recently from a NASA Grand Challenge team award, NASA Cooperative Agreement Number (CAN) NCCS-5-151, through a subcontract from the University of Chicago. We would also like to acknowledge support for computer time from the ASCI program at Livermore and Los Alamos, from the NSF PACI program through NCSA, and also local support from the University of Minnesota's Minnesota Supercomputing Institute. We would also like to acknowledge support to this project for scientific visualization and handling of the very large data sets involved. This support has come from an NSF CISE Research Infrastructure grant, CDA-950297, from the NSF's PACI program through NCSA, from the DoE ASCI program at Livermore (through NCSA), and from the LCSE industrial partners Silicon Graphics, MTI, Ciprico, Seagate Technology, Ancor, and Brocade Communications. We made use of powerful visualization machines at Livermore and at the ACL in Los Alamos together with the LCSE visualization systems to process the data from the large runs performed at these labs. This work was performed in part under the auspices of the U. S. D. o. E. by the Lawrence Livermore National Laboratory under contract NO. W-7405-ENG-48.